

Database on the structure of large ribosomal subunit RNA

Peter De Rijk, Yves Van de Peer, Sabine Chapelle and Rupert De Wachter*

Departement Biochemie, Universiteit Antwerpen (UIA), Universiteitsplein 1, B-2610 Antwerpen, Belgium

ABSTRACT

A database on large ribosomal subunit RNA is made available. It contains 258 sequences. It provides sequence, alignment and secondary structure information in computer-readable formats. Files can be obtained using ftp.

INTRODUCTION

This paper presents a comprehensive database of large ribosomal subunit RNA (further abbreviated as LSU rRNA) structures. Our goal is to offer researchers on-line access to LSU rRNA sequences in the form of an alignment containing secondary structure information, in a format suitable for use in computer programs. Literature references and accession numbers in sequence databases are included as well as taxonomic information.

This database will conceivably be used to perform phylogenetic analysis, to find primers or probes, and to elicit the secondary structure of newly determined sequences. It can also be an invaluable tool to find sequence errors introduced during gel reading or typing. These errors cause anomalies in the sequence or structure which can often be easily detected by alignment to a set of known sequences and comparison of their possible secondary structure.

New entries or updates in the EMBL sequence database (1) are continuously scanned for new LSU rRNA sequences using the Current Sequence Awareness program (a service of the Belgian EMBnet Node). These sequences are used to update older entries, or added to the database as new entries. They are then aligned, and their secondary structure is investigated and incorporated into the alignment using the program DCSE (2). When anomalies or errors are found in the annotations in the sequence libraries, these are corrected when possible. A note indicating the changes made is added.

CONTENTS OF THE DATABASE

Only complete or reasonably complete sequences are being incorporated into the database. Partial sequences are excluded when the combined length of the sequence segments in *Escherichia coli* LSU rRNA homologous to the sequenced segments, amounts to less than 70% of the total *Escherichia coli* sequence. The database currently contains 258 sequences, viz. 42 eukaryal, 16 archaeal, 81 bacterial, 36 plastidial and 83 mitochondrial sequences.

Table 1 shows a list of species for which the LSU rRNA structure is recorded in the database. The same taxonomic

classification is used as in the small ribosomal subunit rRNA database (3). For the domain Eukarya, the taxonomic classification of the species is according to Brusca and Brusca (4) for the Animalia, according to Cronquist (5) for the higher plants, according to Ainsworth *et al.* (6) for the zygomycetes and ascomycetes, according to Moore (7) for the basidiomycetes, and according to Margulis *et al.* (8) for the remaining eukaryotes, viz. the Protoctista.

For the Bacteria and the Archaea, the classification is based on the construction of evolutionary trees, explained into more detail in a previous compilation of small ribosomal subunit RNA sequences (9). In short, evolutionary trees are constructed by the neighbor-joining method (10) for all new sequences retrieved from the EMBL (1) and/or GenBank (11) nucleotide libraries. According to the phylogenetic position of the sequences, they are assigned to one of the taxa listed in Fig. 1 of a previous compilation (9) and described essentially by Woese and coworkers (12, 13).

HETEROGENEITY IN SEQUENCE AND CHAIN LENGTH

Bacterial, archaeal and plastidial LSU sequences have a relatively constant length of approximately 2900 nucleotides. However, eukaryotic sequences show a much greater diversity, ranging in length from sizes comparable to those of the bacteria to over 5000 bases in the *Homo sapiens* sequence. The presence of extra nucleotides seems to be restricted mainly to several extremely variable areas, which occupy a constant position relative to the more conserved parts of the sequences (14,15). Sequence variation is even larger in mitochondria. The molecules found in animal and kinetoplastid mitochondria even miss large parts of the sequence conserved in other LSU rRNAs, and can be under 1000 nucleotides in size. Plant and fungi mitochondrial LSU rRNAs have chain lengths comparable to or larger than those found in bacteria.

SECONDARY STRUCTURE MODEL

Figures 1 and 2 show secondary structure models for a procaryotic (*Escherichia coli*) and a eukaryotic (*Saccharomyces cerevisiae*) LSU rRNA. A core structure is conserved in the majority of eukaryotic and bacterial LSU rRNAs. In the mitochondria of kinetoplastids and animals several helices of this core are absent. Other mitochondria have most of the helices of the core, although the structural variability is higher than among bacteria. The variable insertion regions in Eukaryotes can have

*To whom correspondence should be addressed

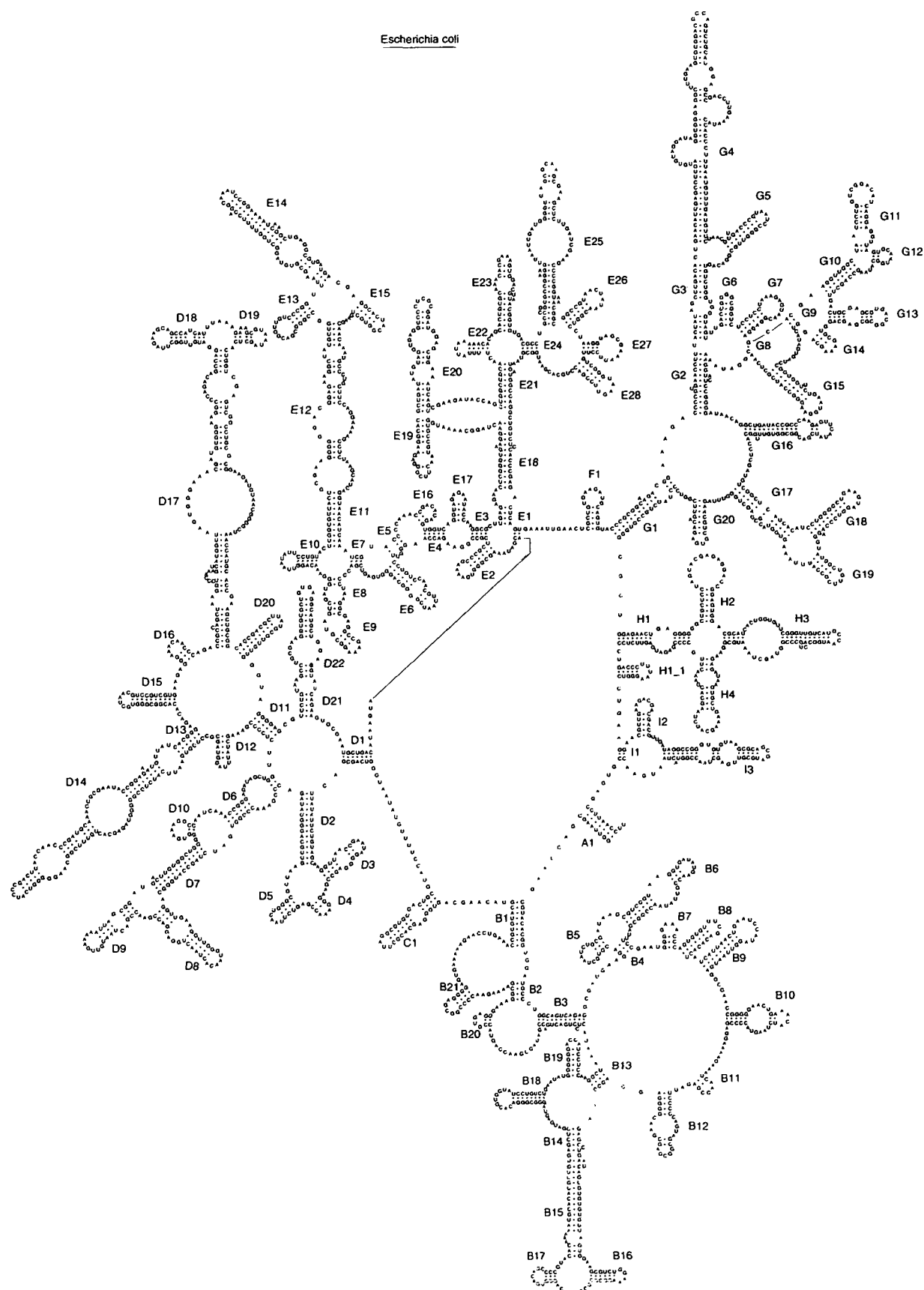


Figure 1. Secondary structure model for *Escherichia coli* LSU rRNA. The sequence is written clockwise from 5' to 3' terminus.

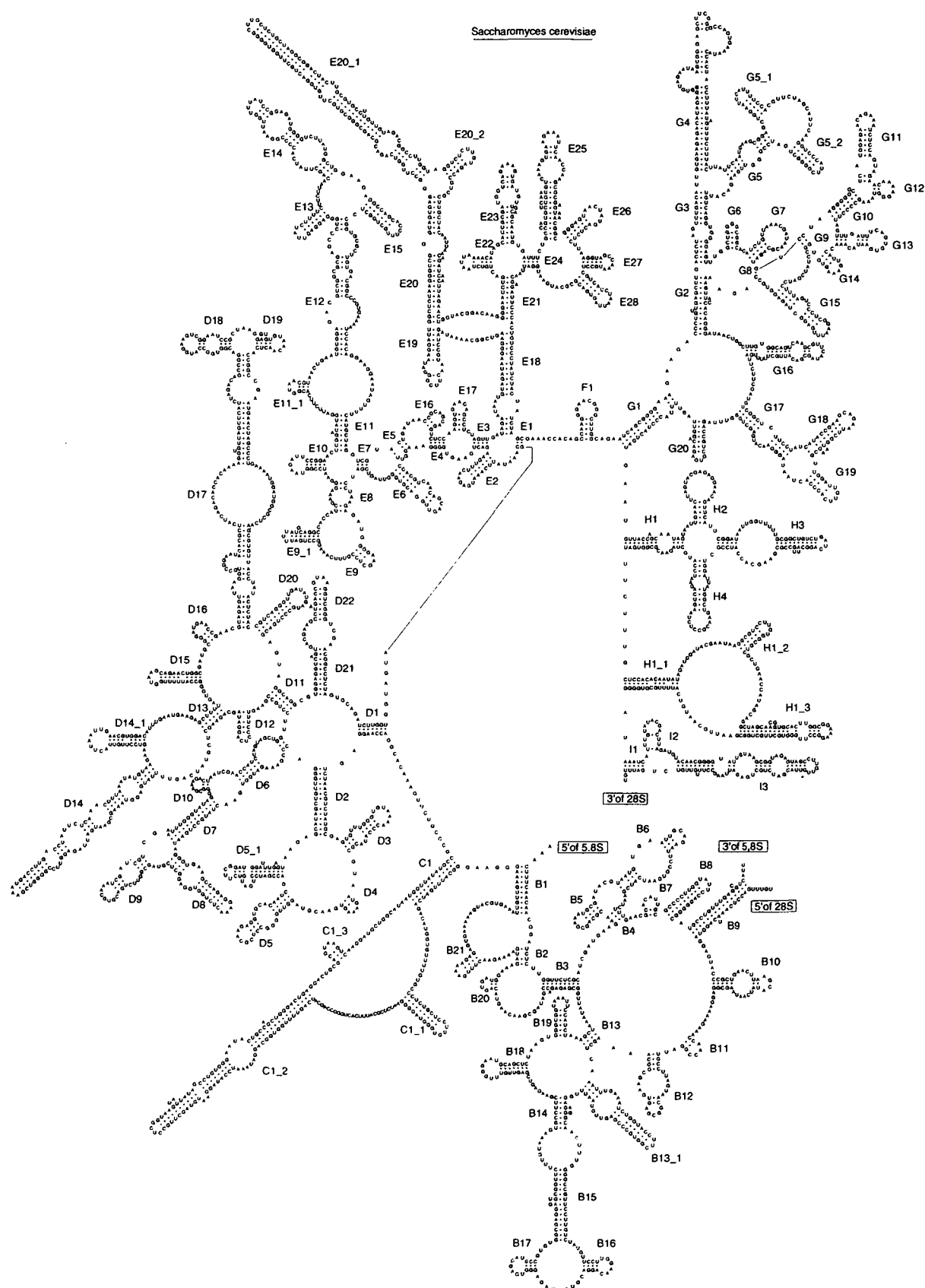


Figure 2. Secondary structure model for *Saccharomyces cerevisiae* LSU rRNA.

Table 1. List of species for which LSU rRNA structure is recorded in the database^a

ARCHAEA	
CRENARCHAEOTA <i>Desulfurococcus mobilis</i> <i>Pyrobaculum islandicum</i> <i>Sulfolobus acidocaldarius</i> <i>Sulfolobus solfataricus</i> <i>Thermophilum pendens</i> DSM 2475 <i>Thermoproteus tenax</i>	METHANOBACTERIALES <i>Methanobacterium thermoautotrophicum</i> METHANOCOCCALES <i>Methanococcus vannielii</i> METHANOMICROBIUM GROUP <i>Methanospirillum hungatei</i> THERMOCOCCALES <i>Thermococcus celer</i> THERMOPLASMA <i>Thermoplasma acidophilum</i>
EURYARCHAEOTA HALOBACTERIA <i>Halobacterium halobium</i> <i>Halobacterium maris-mortui</i> <i>Halococcus morrhuae</i> ATCC 17082 <i>Haloferax volcanii</i> ARCHAEOGLOBALES <i>Archaeoglobus fulgidus</i>	
BACTERIA	
PROTEOBACTERIA ALPHA <i>Bradyrhizobium japonicum</i> DSM 30131 <i>Rhodobacter capsulatus</i> DSM 938 <i>Rhodobacter sphaeroides</i> 1 <i>Rhodobacter sphaeroides</i> 2 <i>Rhodobacter sphaeroides</i> 3 <i>Rhodospseudomonas palustris</i> DSM 126 PROTEOBACTERIA BETA <i>Bordetella avium</i> <i>Bordetella bronchiseptica</i> <i>Bordetella parapertussis</i> <i>Bordetella pertussis</i> <i>Neisseria gonorrhoeae</i> <i>Neisseria meningitidis</i> <i>Pseudomonas cepacia</i> DSM 50181 PROTEOBACTERIA GAMMA <i>Aeromonas hydrophila</i> <i>Escherichia coli</i> 1 <i>Escherichia coli</i> 2 <i>Escherichia coli</i> 3 <i>Escherichia coli</i> 4 <i>Plesiomonas shigelloides</i> NCIMB 9242 <i>Ruminobacter amylophilus</i> <i>Pseudomonas aeruginosa</i> ATCC 10145 <i>Pseudomonas perfectomarina</i> PROTEOBACTERIA EPSILON <i>Campylobacter jejuni</i> ATCC 43431 GRAM POSITIVES AND RELATIVES, HIGH G+C <i>Frankia</i> sp. 1 <i>Frankia</i> sp. 2 <i>Micrococcus luteus</i> <i>Mycobacterium kansasii</i> ATCC 12478 <i>Mycobacterium leprae</i> 1 <i>Mycobacterium leprae</i> 2 <i>Streptomyces ambofaciens</i> ATCC 23877 <i>Streptomyces griseus</i> KCTC 9080 GRAM POSITIVES AND RELATIVES, LOW G+C <i>Bacillus alcalophilus</i> <i>Bacillus anthracis</i> <i>Bacillus cereus</i> NCTC 11143 <i>Bacillus globisporus</i> DSM 4 <i>Bacillus licheniformis</i> DSM 13 <i>Bacillus</i> sp. <i>Bacillus stearothermophilus</i> 1 <i>Bacillus stearothermophilus</i> 2 <i>Bacillus subtilis</i> 1 <i>Bacillus subtilis</i> 2 <i>Clostridium botulinum</i> 1	<i>Clostridium botulinum</i> 2 <i>Clostridium botulinum</i> 3 ATCC 25765 <i>Clostridium botulinum</i> 4 <i>Clostridium botulinum</i> 5 NCTC 7272 <i>Clostridium tyrobutyricum</i> <i>Lactobacillus confusus</i> NCDO 1586 <i>Lactobacillus delbrueckii</i> subsp. <i>bulgari</i> <i>Lactococcus lactis</i> subsp. <i>cremoris</i> DSM 20069 <i>Lactococcus lactis</i> subsp. <i>lactis</i> 1 DSM 20481 <i>Lactococcus lactis</i> subsp. <i>lactis</i> 2 <i>Leuconostoc carnosum</i> <i>Leuconostoc mesenteroides</i> <i>Leuconostoc oenos</i> <i>Leuconostoc paramesenteroides</i> <i>Listeria monocytogenes</i> 1 ATCC 19115 <i>Listeria monocytogenes</i> 2 <i>Mycoplasma flocculare</i> <i>Mycoplasma hyopneumoniae</i> ATCC 27719 <i>Mycoplasma pneumoniae</i> <i>Pectinatus frisingensis</i> DSM 20465 <i>Peptococcus niger</i> DSM 20475 <i>Staphylococcus aureus</i> ATCC 12208 <i>Staphylococcus carnosus</i> DSM 20501 <i>Streptococcus oralis</i> DSM 20066 <i>Streptococcus parauberis</i> NCDO 2020 <i>Streptococcus thermophilus</i> DSM 20617 <i>Streptococcus uberis</i> NCDO 2038 CYANOBACTERIA <i>Anacystis nidulans</i> FLAVOBACTERIA AND RELATIVES <i>Flavobacterium odoratum</i> ATCC 4651 <i>Flexibacter flexilis</i> ATCC 23079 GREEN SULFUR <i>Chlorobium limicola</i> ATCC 8327 PLANCTOMYCES AND RELATIVES <i>Pirellula marina</i> SPIROCHETES <i>Borrelia burgdorferi</i> 1 <i>Borrelia burgdorferi</i> 2 <i>Borrelia burgdorferi</i> 3 <i>Leptospira interrogans</i> RADIORESISTANT MICROCOCCI AND RELATIVES <i>Thermus thermophilus</i> THERMOTOBALES <i>Thermotoga maritima</i>
EUKARYA	
ANIMALIA CHORDATA VERTEBRATA	MAGNOLIOPSIDA <i>Arabidopsis thaliana</i> <i>Brassica napus</i>

Table 1. (cont.)

AMPHIBIA <i>Xenopus borealis</i> <i>Xenopus laevis</i> 1 <i>Xenopus laevis</i> 2	<i>Citrus limon</i> <i>Fragaria ananassa</i> <i>Lycopersicon esculentum</i> <i>Sinapis alba</i>
MAMMALIA <i>Homo sapiens</i> <i>Mus musculus</i> <i>Rattus norvegicus</i>	PROTOCTISTA APICOMPLEXA COCCIDIA <i>Toxoplasma gondii</i> 1 <i>Toxoplasma gondii</i> 2 <i>Toxoplasma gondii</i> 3
UROCHORDATA <i>Herdmania momus</i>	CHYTRIDIOMYCOTA OOMYCOTA <i>Phytophthora megasperma</i>
ARTHROPODA INSECTA <i>Aedes albopictus</i> <i>Drosophila melanogaster</i>	CILIOPHORA <i>Tetrahymena pyriformis</i> <i>Tetrahymena thermophila</i>
NEMATODA SECERNENTEA <i>Caenorhabditis elegans</i>	DICTYOSTELIDA <i>Dictyostelium discoideum</i>
FUNGI ZYGYMYCOTINA ZYGYMYCETES <i>Mucor racemosus</i>	DINOFLAGELLATA <i>Prorocentrum micans</i>
ASCOMYCOTINA HEMIASCOMYCETES <i>Candida albicans</i> <i>Saccharomyces cerevisiae</i> <i>Schizosaccharomyces pombe</i> 1 <i>Schizosaccharomyces pombe</i> 2	EUGLENIDA <i>Euglena gracilis</i>
UNCERTAIN AFFILIATION <i>Pneumocystis carinii</i>	PLASMODIAL SLIME MOLDS MYXOMYCOTA <i>Didymium iridis</i> <i>Physarum polycephalum</i>
BASIDIOMYCOTINA HETEROBASIDIOMYCETES <i>Cryptococcus neoformans</i> 1 <i>Cryptococcus neoformans</i> 2	RHIZOPODA LOBOSEA <i>Entamoeba histolytica</i>
PLANTAE MAGNOLIOPHYTA LILIOPSIDA <i>Oryza sativa</i>	ZOOMASTIGINA DIPLOMONADIDA <i>Giardia ardeae</i> <i>Giardia intestinalis</i> <i>Giardia muris</i>
	KINETOPLASTIDA <i>Critidia fasciculata</i> <i>Trypanosoma brucei</i>
PLASTIDS	
PLANTAE BRYOPHYTA MARCHANTIOPSIDA <i>Marchantia polymorpha</i>	<i>Chlamydomonas humicola</i> <i>Chlamydomonas indica</i> <i>Chlamydomonas iyengarii</i> <i>Chlamydomonas komma</i> <i>Chlamydomonas mexicana</i> <i>Chlamydomonas moewusii</i> <i>Chlamydomonas pallidostigmatica</i> <i>Chlamydomonas pteridii</i> <i>Chlamydomonas pilschuanii</i> <i>Chlamydomonas reinhardtii</i> <i>Chlamydomonas sp.</i> <i>Chlamydomonas starrii</i> <i>Chlamydomonas zebra</i> <i>Chlorella ellipsoidea</i>
MAGNOLIOPHYTA LILIOPSIDA <i>Oryza sativa</i> <i>Zea mays</i>	PHAEOPHYTA <i>Pylaiella littoralis</i>
MAGNOLIOPSIDA <i>Alnus incana</i> <i>Conopholis americana</i> <i>Epifagus virginiana</i> 1 <i>Epifagus virginiana</i> 2 <i>Nicotiana tabacum</i> 1 <i>Nicotiana tabacum</i> 2 <i>Pisum sativum</i>	EUGLENIDA <i>Astasia longa</i> <i>Euglena gracilis</i> 1 <i>Euglena gracilis</i> 2 <i>Euglena gracilis</i> 3 <i>Euglena gracilis</i> 4
PROTOCTISTA CHLOROPHYTA CHLOROPHYCEAE <i>Nanoarchaeum eucaryotum</i> <i>Chlamydomonas eugametos</i> <i>Chlamydomonas frankii</i> <i>Chlamydomonas gottliebi</i> <i>Chlamydomonas gelatinosa</i>	RHODOPHYTA <i>Palmaria palmata</i>
MITOCHONDRIA	
ANIMALIA CHORDATA VERTEBRATA	MOLLUSCA BIVALVIA

Table 1. (cont.)

MAMMALIA	<i>Mytilus edulis</i>
<i>Aepyceros melampus</i>	NEMATODA
<i>Antilocapra americana</i>	SECERNENTEA
<i>Balaenoptera musculus</i>	<i>Ascaris suum</i>
<i>Balaenoptera physalus</i>	<i>Caenorhabditis elegans</i>
<i>Bos taurus</i>	
<i>Boselaphus tragocamelus</i>	FUNGI
<i>Capra hircus</i>	ASCOMYCOTINA
<i>Cephalophus maxwelli</i>	HEMIASCOMYCETES
<i>Cervus unicolor</i>	<i>Saccharomyces cerevisiae</i> 1
<i>Damalisus dorcas</i>	<i>Saccharomyces cerevisiae</i> 2
<i>Didelphis virginiana</i>	<i>Schizosaccharomyces pombe</i>
<i>Gazella thomsoni</i>	
<i>Halichoerus grypus</i>	PLECTOMYCETES
<i>Homo sapiens</i> 1	<i>Aspergillus nidulans</i>
<i>Homo sapiens</i> 2	<i>Penicillium chrysogenum</i>
<i>Homo sapiens</i> 3	
<i>Hydropotes inermis</i>	PYRENOMYCETES
<i>Kobus ellipsiprymnus</i>	<i>Neurospora crassa</i>
<i>Madoqua kirkii</i>	<i>Podospira anserina</i>
<i>Muntiacus reevesi</i>	
<i>Mus musculus</i>	PLANTAE
<i>Odocoileus virginianus</i>	BRYOPHYTA
<i>Phoca vitulina</i>	MARCHANTIOPSIDA
<i>Rattus norvegicus</i> 1	<i>Marchantia polymorpha</i>
<i>Rattus norvegicus</i> 2	
<i>Tragelaphus imberbis</i>	MAGNOLIOPHYTA
<i>Tragulus napu</i>	LILIOPSIDA
	<i>Triticum aestivum</i>
AVES	<i>Zea mays</i>
<i>Anas platyrhynchos</i>	
<i>Cairina moschata</i>	MAGNOLIOPSIDA
<i>Gallus gallus</i>	<i>Oenothera berteriana</i>
AMPHIBIA	PROTOCTISTA
<i>Rana catesbeiana</i>	APICOMPLEXA
<i>Xenopus laevis</i>	HEMATOZOA
	<i>Plasmodium falciparum</i>
OSTEICHTHYES	
<i>Crossostoma lacustre</i>	CHLOROPHYTA
<i>Cyprinus carpio</i>	CHLOROPHYCEAE
<i>Neoceratodus forsteri</i> 1	<i>Chlamydomonas eugametos</i>
<i>Neoceratodus forsteri</i> 2	<i>Chlamydomonas reinhardtii</i>
<i>Protopterus</i> sp.	<i>Prototheca wickerhamii</i>
<i>Latimeria chalumnae</i>	<i>Scenedesmus obliquus</i>
ECHINODERMATA	CILIOPHORA
ECHINOIDEA	<i>Paramecium aurelia</i>
<i>Paracentrotus lividus</i>	<i>Paramecium primaurelia</i> 1
<i>Strongylocentrotus purpuratus</i>	<i>Paramecium tetraurelia</i> 2
	<i>Tetrahymena pyriformis</i> 1
ARTHROPODA	<i>Tetrahymena pyriformis</i> 2
MALACOSTRACA	
<i>Artemia franciscana</i>	ZOOMASTIGINA
<i>Artemia salina</i>	KINETOPLASTIDA
	<i>Crithidia fasciculata</i>
INSECTA	<i>Crithidia oncopelti</i>
<i>Aedes albopictus</i>	<i>Herpetomonas mariadeanei</i>
<i>Apis mellifera</i>	<i>Herpetomonas megaseliae</i>
<i>Apis mellifera ligustica</i>	<i>Herpetomonas muscarum</i>
<i>Drosophila melanogaster</i>	<i>Herpetomonas samuelpessoai</i>
<i>Drosophila yakuba</i>	<i>Leishmania tarentolae</i>
<i>Locusta migratoria</i>	<i>Leptomonas</i> sp.
<i>Spodoptera frugiperda</i>	<i>Trypanosoma brucei</i> 1
	<i>Trypanosoma brucei</i> 2

*In some cases, species names are listed several times followed by a sequential number, because multiple LSU rRNA sequences have been determined, usually by different authors. These sequences are not necessarily the same because they may originate from different varieties or strains, or from different genes, of the same species. The taxonomy followed for the three domains Eukarya, Archaea, and Bacteria, is as explained in the text. Plastidial and mitochondrial structures are listed according to the systematics followed for the host organism. In the case of Archaea and Bacteria, the species name is followed by the culture collection name and number if specified by the author.

differences in length of up to 900 bases. The structure of some of these regions has not yet been conclusively determined. The alignment and proposed secondary structure of the mitochondrial LSU rRNAs is less dependable because of the larger variability in both length and sequence.

The secondary structure of the molecule is treelike, with the helices forming branches which end either in a hairpin or in a multibranched loop. The stem of the tree joins the 5' and 3' end of bacterial LSU rRNAs. From this stem emanates a central multibranched loop. In Eukarya, and probably in Archaea the

stem helix is not present, but the central loop is. The following provisional helix numbering system is used in Figs 1 and 2. Structures branching from the central loop are labeled A to I, starting with the stem helix. Within each of these structures, helices bear a different number when they are separated by a multibranched loop. All numbering is sequentially from 5' to 3'. Structural elements specific to certain taxa are named after the preceding core helix followed by an underscore and number. The helix numbering may have to be revised if additional structural elements are identified in the future.

AVAILABILITY AND FORMAT OF THE DATABASE

The LSU rRNA database will be made available through anonymous ftp on the server uiam3.uia.ac.be (143.169.8.1). The files will also be made available to the EMBL nucleotide library for distribution. Researchers who cannot obtain the database through these channels, can request the database or parts thereof on magnetic media from the authors. The authors can be contacted by electronic mail to dwachter@reks.uia.ac.be or rrna@reks.uia.ac.be. On the server, a file called 'readme' will be present which describes the latest state of the database, giving the contents of the files and directories, and a description of the programs available for format conversion, alignment editing (2) and phylogenetic tree construction (16).

In order to simplify access to the database, each sequence is stored in a separate file, together with information about this sequence. The names of these files are produced from the species name by taking characters of the genus and species names. These are preceded by a code describing the phylogenetic group to which the species belongs. This makes it possible to either retrieve specific sequences using the full file name, or to retrieve a set of sequences belonging to a phylogenetic group using wild cards. Several sequence files can be integrated into one alignment using a program available on the server.

The format of the files is very simple, so that the files can be used readily by computer programs, or can easily be converted to formats used by specific programs. The files start with a few header lines which contain data about the sequence such as the accession number and literature reference. These are followed by the organism name. The sequence comes next. It consists of a range of nucleotide symbols interspersed with gap symbols necessary for alignment. The sequence end is indicated by an asterisk. The beginning and end of secondary structure elements are indicated by insertion of special symbols. Special 'helix numbering' files are present for researchers who wish to use the secondary structure information. When these are incorporated into an alignment, they indicate the name of each different helix segment.

When a sequence consists of several fragments resulting from processing, or of several exons, the sequence of each part ends with an asterisk, and has its own header containing the accession number, literature reference and a description of the sequence segment. However, the segments are stored in the same file and have the same organism name.

Users of the database are requested to cite this paper.

ACKNOWLEDGEMENTS

Peter De Rijk is research assistant of the National Fund for Scientific Research. Our research was supported by the Programme on Interuniversity Poles of Attraction (contract 23)

of the Office for Science Policy Programming of the Belgian State, and by the Fund for Collective Fundamental Research.

REFERENCES

1. Rice, C.M., Fuchs, R., Higgins, D.G., Stoehr, P.J. and Cameron, G.N. (1993) *Nucleic Acids Res.* **21**, 2967–2971.
2. De Rijk P. and De Wachter, R. (1993) *Comput. Applic. Biosci.*, **9**, 735–740.
3. Van de Peer, Y., Van den Broeck, I., De Rijk, P. and De Wachter, R. *Nucleic Acids Res.*, this issue
4. Brusca, R.C. and Brusca, G.J. (1990) *Invertebrates*, Sinauer Associates, Inc. Sunderland.
5. Cronquist, A. (1971) *Introductory Botany*, Harper & Row, New York.
6. Ainsworth, G.C., Sparrow, F.K. and Sussman, A.S. (1973), *The Fungi: an Advanced Treatise*, Academic Press, New York, Vol. 4A.
7. Moore, R.T. (1988) in Moriarty, Ch. (ed.), *Taxonomy putting plants and animals in their place*. Royal Irish Academy, Dublin, pp. 61–88.
8. Margulis, L., Corliss, J.O., Melkonian, M. and Chapman, D.J. (eds.) (1990) *Handbook of Protozoa*, Jones and Bartlett Publishers, Boston.
9. Neefs, J.-M., Van de Peer, Y., De Rijk, P., Chapelle, S. and De Wachter, R. (1993) *Nucleic Acids Res.* **21**, 3025–3049.
10. Saitou, N. and Nei, M. (1987) *Mol. Biol. Evol.* **4**, 406–425.
11. Benson, D., Lipman, D.J. and Ostell, J. (1993) *Nucleic Acids Res.* **21**, 2963–2965.
12. Woese, C.R. (1987) *Microbiol. Rev.* **51**, 221–271.
13. Olsen, G.J., Woese, C.R. and Overbeek, R. (1994) *J. Bacteriol.* **176**, 1–6.
14. Veldman, G.M., Klootwijk, J., de Regt, V.C.H.F., Planta, R.J., Brantlant, C., Krol, A. and Ebel, J.-P. (1981) *Nucleic Acids Res.*, **9**, 6935–6952.
15. Michot, B., Hassouna, N. and Bachellerie, J.-P. (1984) *Nucleic Acids Res.*, **12**, 4259–4279.
16. Van de Peer, Y. and De Wachter, R. (1993) *Comput. Applic. Biosci.*, **9**, 177–182.